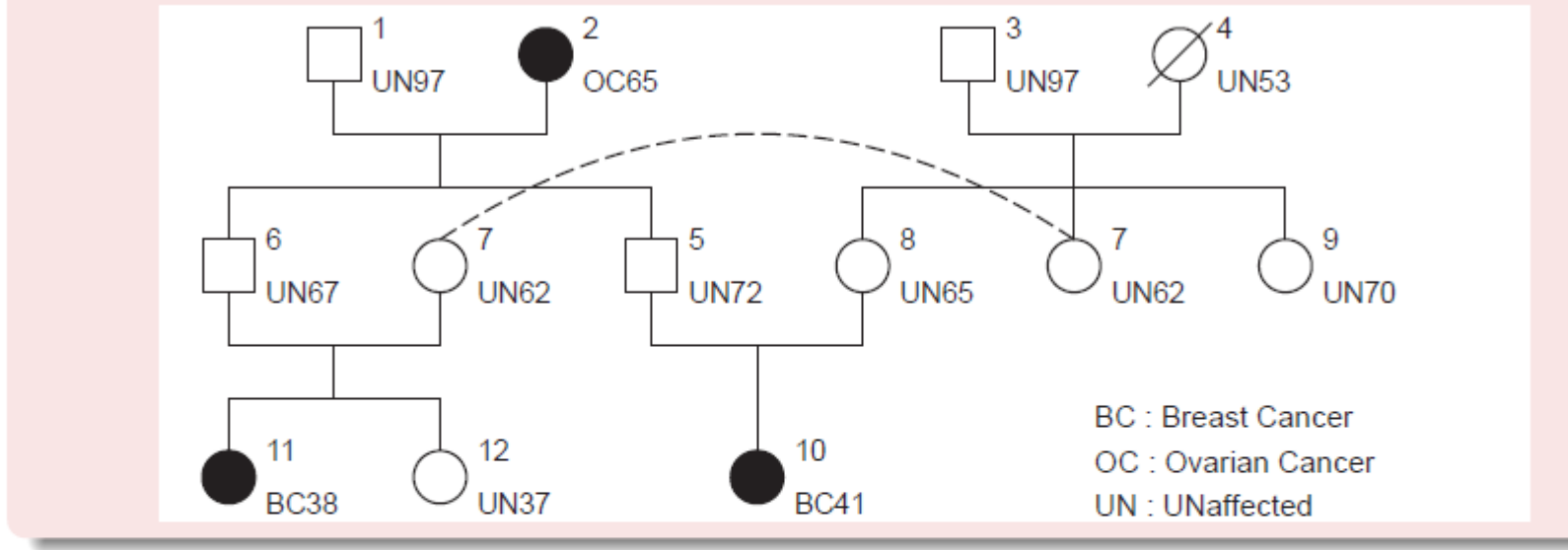


Context of the breast cancer

- 1st cancer in women. Affects 54,000 women in France each year
- Complex disease due to an accumulation of mutations in oncogenes and/or suppressor of tumor genes (BRCA 1/2, PALB2, RAD51, etc.)
- Inherited mutation** in 10 to 15% of the cases leading more often to severe family histories (FH)

Data structure in genetic diseases modeling

Family History (FH) : Pedigree + Sequencing (seq) + Survival (Y)



- #### Variant data (var)
- Known variant or VUS (Uncertain Significance)
 - Database
 - Functional tests
 - Molecular dynamics

- #### Cancer pathology (patho)
- MSI status (Lynch) (MicroSatellite Instability)
 - BRAF mutation (Lynch)
 - ER, PR, HER2 status (Breast)
 - invasive/in situ (Breast)

V : the set of variant status

var : the set of variant data

X : the set of all possible values for the genotypes

seq : the set of sequencing data

Y : the set of phenotypes (survival data) } $FH = \{seq, Y\}$

$patho$: the set of pathology reports

$$\mathbb{P}(var, V, X, seq, Y, patho) = \prod_j \mathbb{P}(var_j | V_j) \mathbb{P}(V_j) \times \prod_{i, Y_i \neq UN} \mathbb{P}(patho_i | X_i) \times \prod_i \mathbb{P}(X_i | X_{pat_i}, X_{mat_i}) \mathbb{P}(seq_i | X_i) \mathbb{P}(Y_i | X_i)$$

compute \Rightarrow posterior variant status carrier risk tumoral risk

The Claus-Easton model

- Developed by Claus et al. (1991) and Easton et al. (1993)
- Used in first intention at the Curie Institute
- Focuses on the genotypes (X) and the phenotypes (Y)
- Autosomal, biallelic, dominant mode of inheritance
- Estimated allele frequency ($f = 0.33\%$); hazard functions per genotype (λ_0 and λ_1) derived from Easton's estimated densities

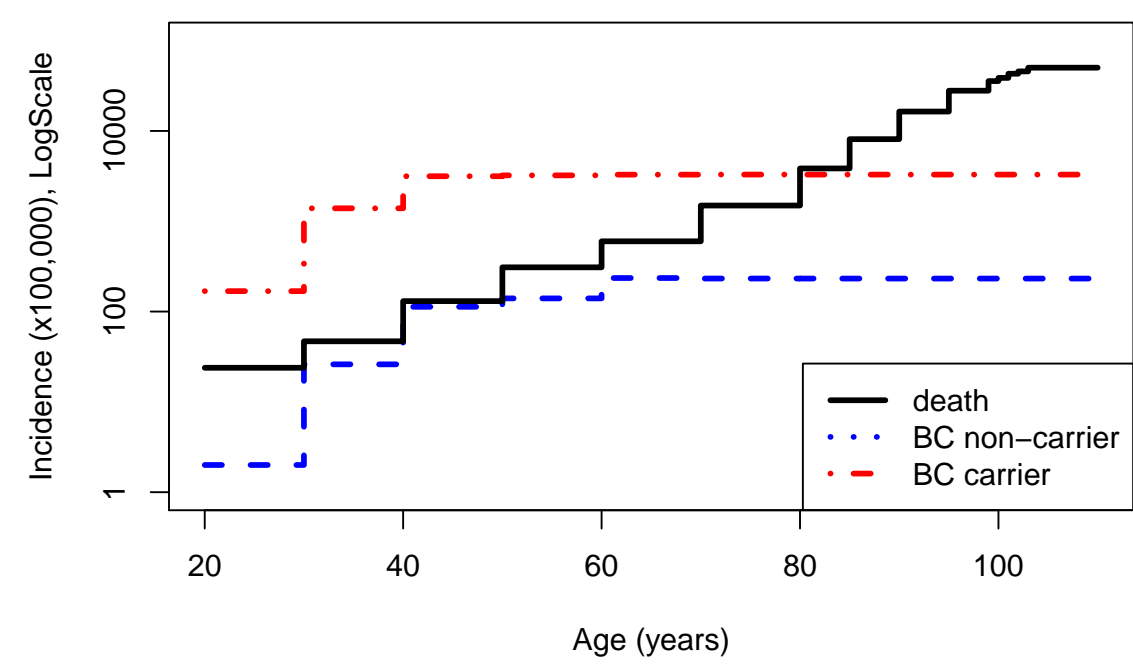


Figure 1: Annual death incidences in the French female population (INED, 2017) and annual breast cancer incidences for non-carriers and carriers estimated from Easton et al. (1993).

Our objectives: Implement the Claus-Easton model in a Bayesian network (sum-product algorithm) combined with survival data and develop a user-friendly interface.

Implementation

$$\mathbb{P}(X, Y) = \prod_i \mathbb{P}(X_i | X_{pat_i}, X_{mat_i}) \mathbb{P}(Y_i | X_i)$$

The genotypes (unobserved): $X = (X_i)_{i=1, \dots, n} \in \{00, 01, 10, 11\}^n$

- Mode of inheritance and allele frequency from the Claus-Easton literature
- Hardy-Weinberg** for the founders (assumption)
- Mendelian transmission** for the offsprings (assumption)

The phenotypes (observed) (FH)

$$Y_i(\text{survival data}) = \begin{cases} \{T_i > \tau_i\} & \text{if } i \text{ is censored (UN) at age } \tau_i \\ \{T_i = \tau_i\} & \text{if } i \text{ is affected (BC or OC) at age } \tau_i \end{cases}$$

with T_i being the age at disease onset for individual i .

Incidence

$$\lambda(t) = \begin{cases} \lambda_0(t) & \text{for } X = 00 \text{ (non-carrier (NC))} \\ \lambda_1(t) & \text{for } X \neq 00 \text{ (carrier (C))} \end{cases}$$

Survival functions of T_i :

$$\begin{cases} S_0(t) = \exp(-\int_0^t \lambda_0(u) dt) & \text{for NC} \\ S_1(t) = \exp(-\int_0^t \lambda_1(u) dt) & \text{for C} \end{cases}$$

Conditional probabilities:

- For a **censored** individual at age τ_i : $\mathbb{P}(Y_i | X_i) = \begin{cases} S_0(\tau_i) & \text{for NC} \\ S_1(\tau_i) & \text{for C} \end{cases}$
- For an **affected** individual at age τ_i : $\mathbb{P}(Y_i | X_i) = \begin{cases} S_0(\tau_i) \lambda_0(\tau_i) & \text{for NC} \\ S_1(\tau_i) \lambda_1(\tau_i) & \text{for C} \end{cases}$

Computation of the posterior carrier probability

Problematic

We denote by $K_i(X_i, X_{pat_i}, X_{mat_i})$, the potential related to i .
 $K_i(X_i, X_{pat_i}, X_{mat_i}) = \sum_{Y_i \in FH} \mathbb{P}(X_i | X_{pat_i}, X_{mat_i}) \mathbb{P}(Y_i | X_i)$

Using the **Bayes rule**, for any individual j ,

$$\mathbb{P}(X_j = x_j | FH) = \frac{\mathbb{P}(X_j = x_j, FH)}{\mathbb{P}(FH)} = \frac{\sum_{X \setminus X_j} \prod_i K_i(X_i, X_{pat_i}, X_{mat_i})}{\sum_X \prod_i K_i(X_i, X_{pat_i}, X_{mat_i})}$$

With $X \in \{00, 10, 01, 11\}^n \rightarrow 4^n$ configurations

The sum-product algorithm (Koller and Friedman, 2009) is equivalent to the latest version of Elston-Stewart algorithm (Totir et al., 2009)

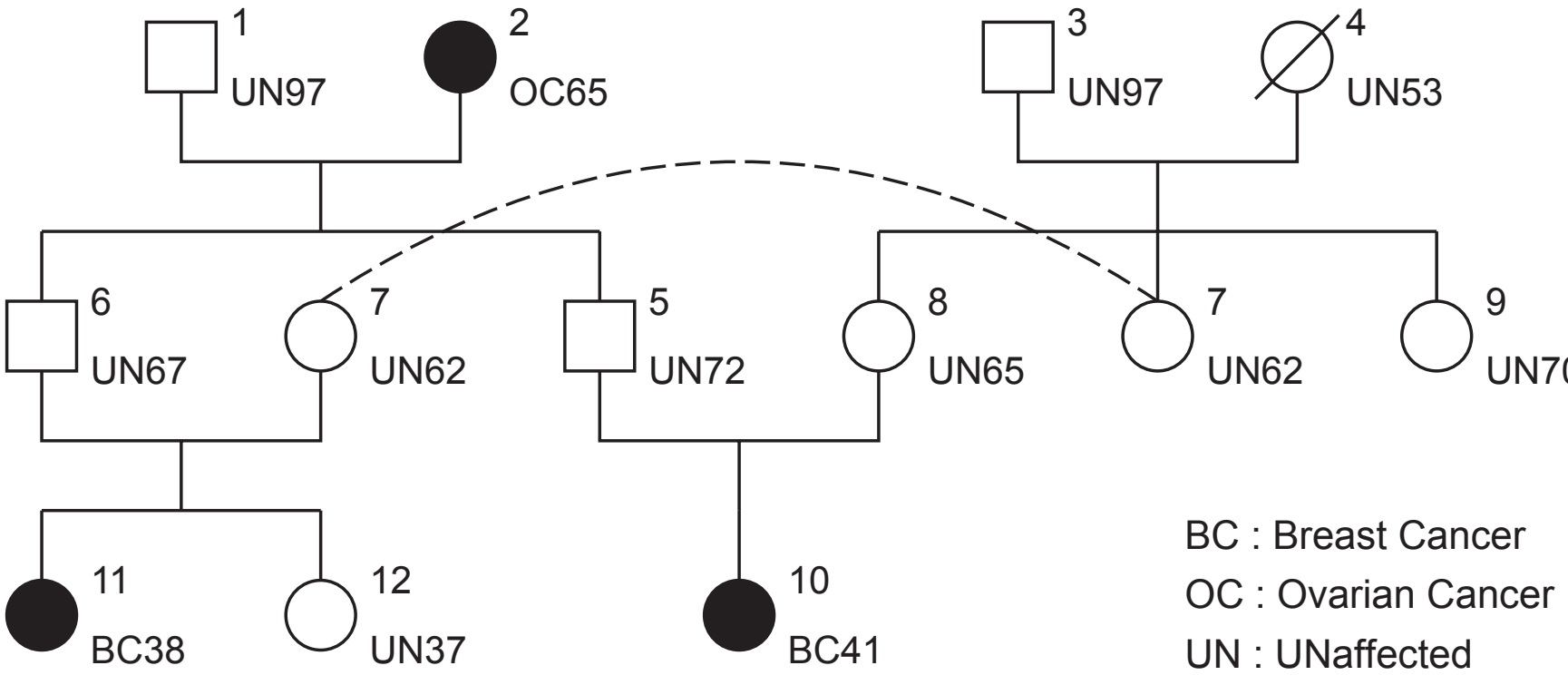
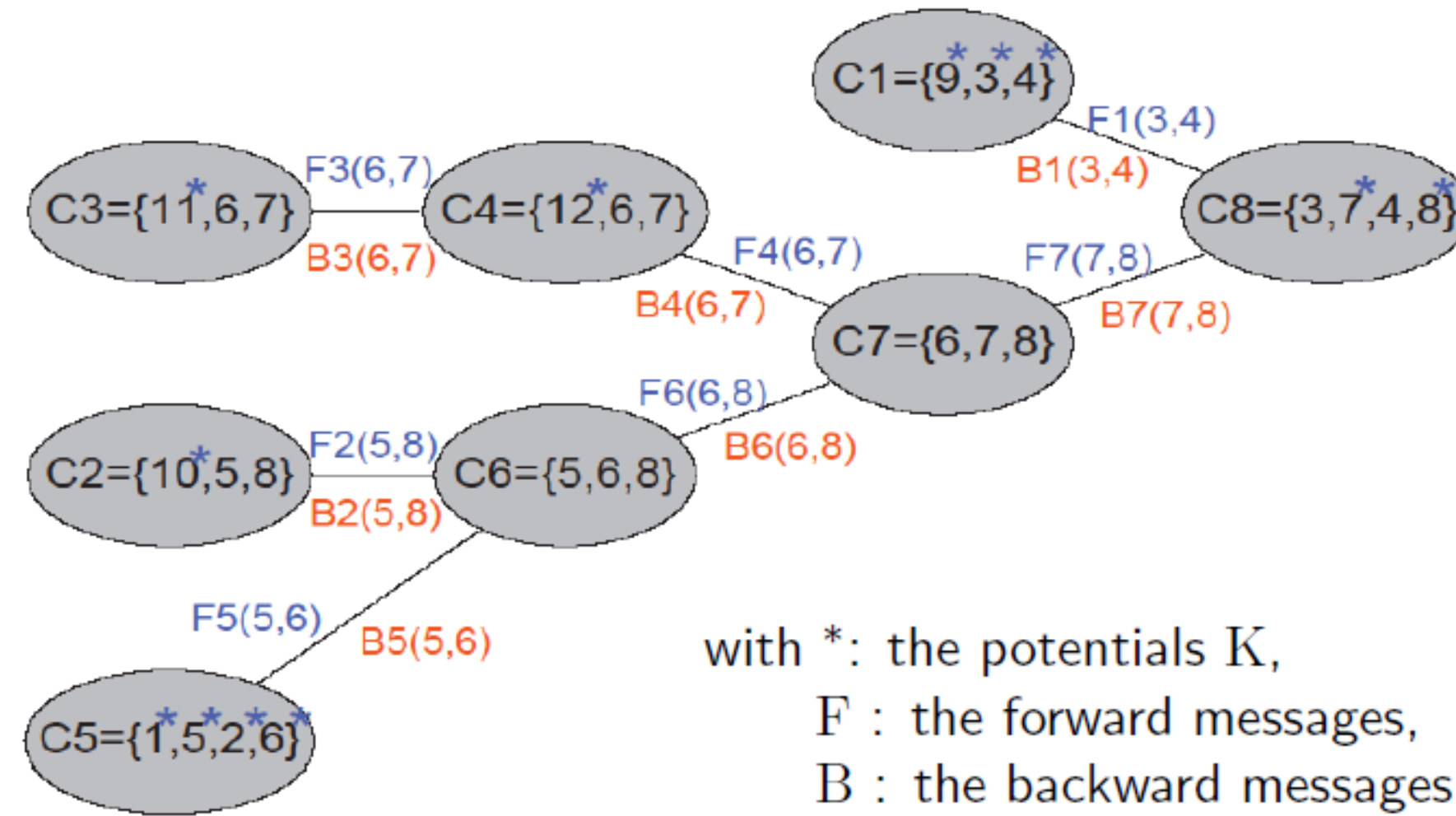


Figure 2: Hypothetical family with a severe FH.

Using the conditional independencies between the individuals and the minimum fill-in heuristic, we obtain the following junction tree:



with *: the potentials K ,
F: the forward messages,
B: the backward messages.

Figure 3: One junction-tree obtained from our hypothetical family.

e.g. $F_3(6, 7) = \sum_{X_{11}} K_{11}(X_{11}, X_6, X_7)$
 $F_4(6, 7) = \sum_{X_{12}} K_{12}(X_{12}, X_6, X_7) F_3(6, 7)$ $B_4(6, 7) = \sum_{X_6} F_6(6, 8) B_7(7, 8)$

All F&B messages computed once for any later marginal or joint distribution needed \Rightarrow Complexity $\mathcal{O}(4^n) \rightarrow \mathcal{O}(n \times 4^k)$
k: tree-width = 3 in general (2 parents & 1 child in a clique),
4 in case of a loop (mating loop, consanguinity),
5 or more exceptionally (several loops).

e.g.: $\mathbb{P}(X_6, X_7, X_8 | FH) \propto F_4(6, 7) F_6(6, 8) B_7(7, 8)$
 $\mathbb{P}(X_6 | FH) \propto \sum_{X_8} F_6(6, 8) B_6(6, 8)$

Results: The interface

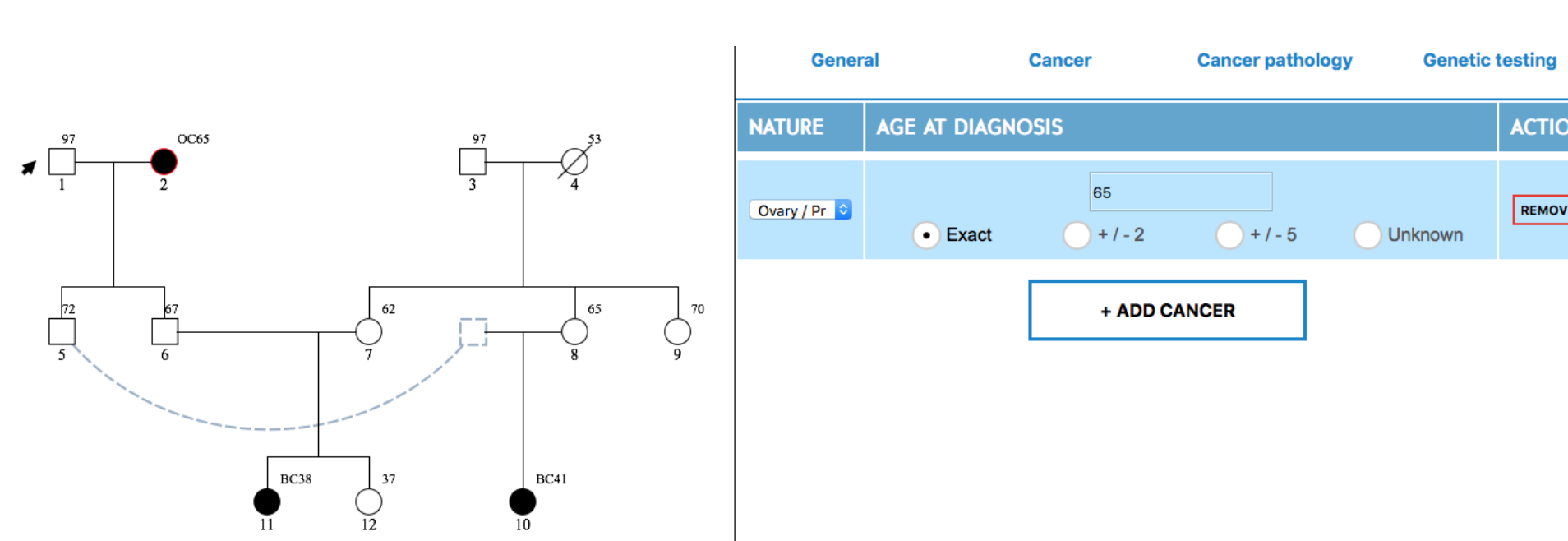


Figure 4: Example of a data entry of a hypothetical family history with the interactive interface.

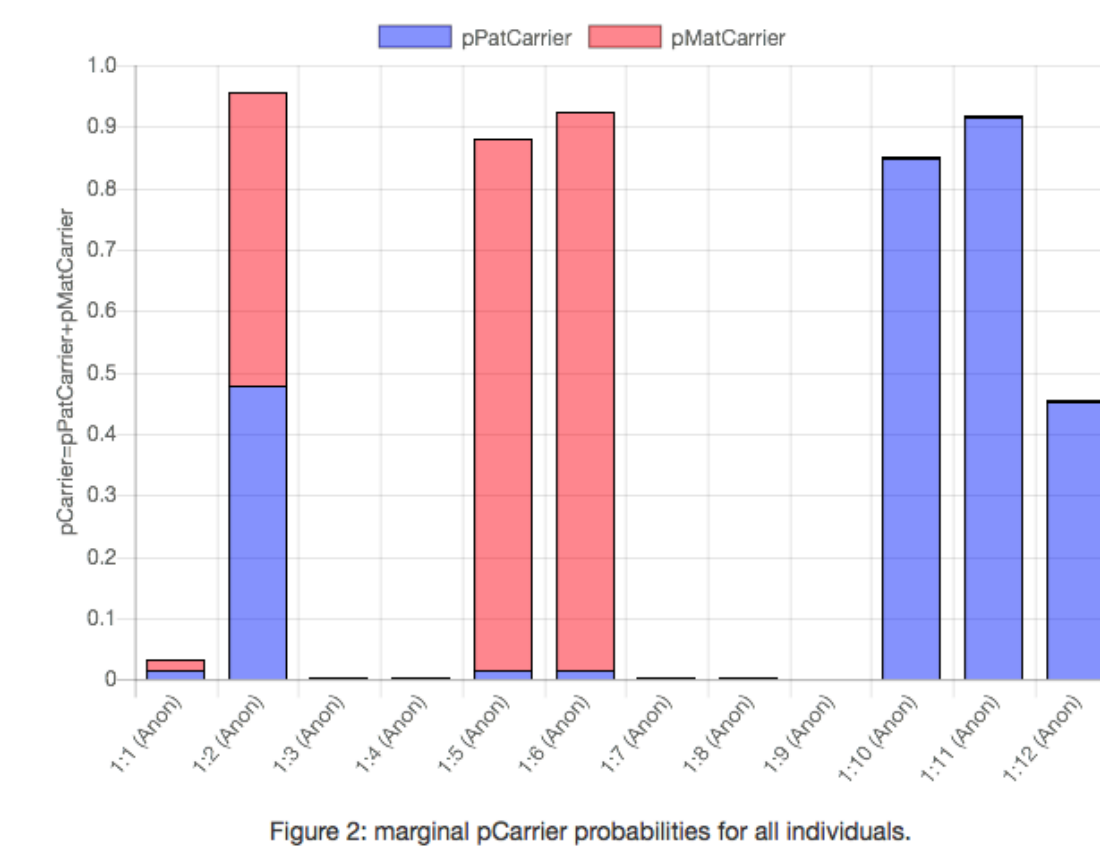


Figure 2: marginal pCarrier probabilities for all individuals.

Ind	pPatCarrier	pMatCarrier	pCarrier
1:1 (Anon) M 97	0.0160167	0.0160167	0.0320333
1:2 (Anon) F OC65	0.47883	0.47883	0.95766
1:3 (Anon) M 97	0.00258042	0.00258042	0.00516084
1:4 (Anon) F 53	0.00142826	0.00142826	0.00285651
1:5 (Anon) M 72	0.0153911	0.865806	0.881197
1:6 (Anon) M 67	0.0160017	0.908384	0.924385
1:7 (Anon) F 62	0.00388557	0.00170409	0.00478965

Figure 5: Probability of being a carrier in a barplot (on top) and in a table (at the bottom) obtained with the interface for each individual of our hypothetical family.

Results - Carrier predisposition

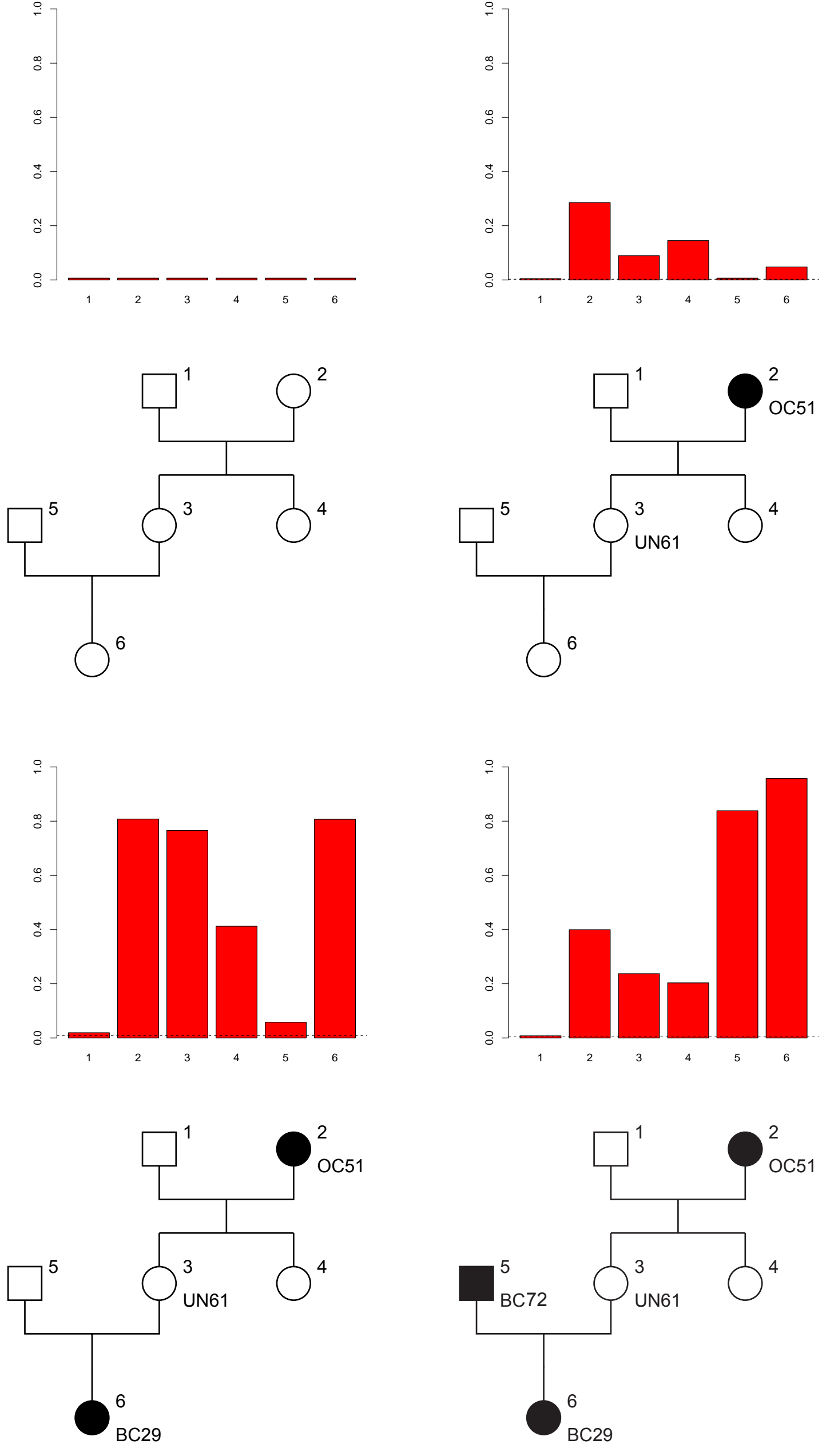


Figure 6: Probability of being a carrier for each individual of a family with evolving FH.

Computation of the disease risks prediction

We denote by $\pi(\tau) = \mathbb{P}(X_i \neq 00(\text{carrier}) | FH)$

- With no competing risk of death
 $\mathbb{P}(T \leq t | FH) = 1 - S(t | FH)$ with
 $S(t | FH) = \sum_{X_i} \mathbb{P}(T > t, X_i | FH) = \pi(\tau) \frac{S_0(t)}{S_0(\tau)} + (1 - \pi(\tau)) \frac{S_1(t)}{S_1(\tau)}$
- Quantities needed to compute the competing risk
 $\pi(t | FH) = \frac{\pi(\tau) S_1(t)}{S(t | FH) S_1(\tau)}$
 $\lambda_{\text{disease}}(t | FH) = \pi(t | FH) \lambda_1(t) + (1 - \pi(t | FH)) \lambda_0(t)$
- With competing risk of death
 $T^* = \min(T_{\text{disease}}, T_{\text{death}})$
 $\lambda_{\text{both}}(t | FH) = \lambda_{\text{disease}}(t | FH) + \lambda_{\text{death}}(t)$
 $\mathbb{P}(T \leq t | FH) = \int_0^t S_{\text{both}}(u) \lambda_{\text{disease}}(u) du$

Results - Tumoral risk

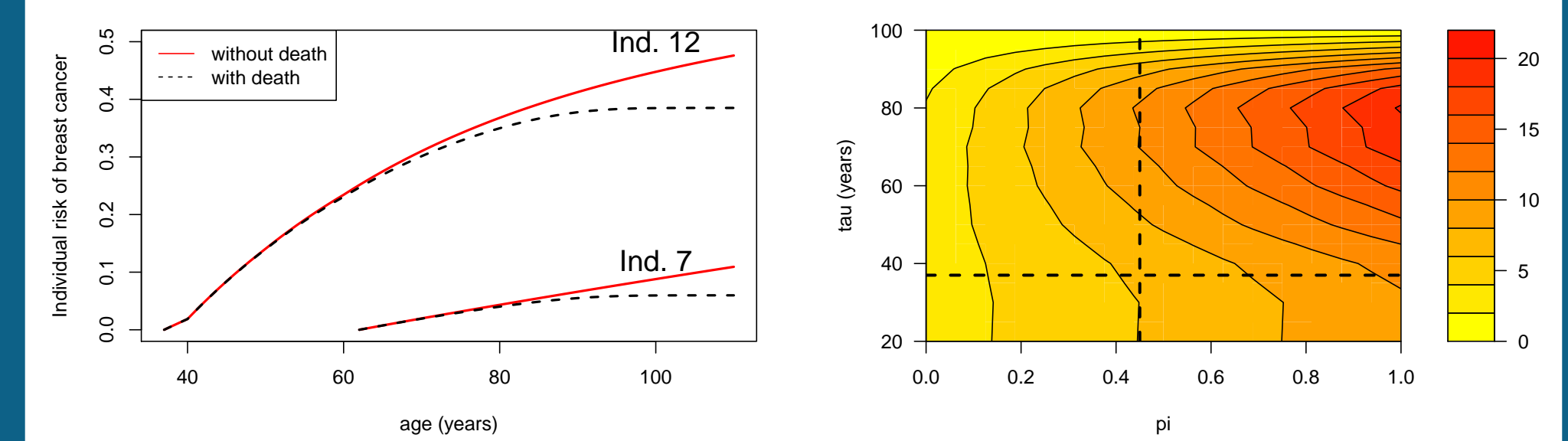


Figure 7: Left-panel: Disease risk for individuals 12 and 7 in our hypothetical family. The risk is computed with and without the consideration of the competing risk of death. Right-panel: Percentage of error made (difference) while computing the disease risk up to the age 100 without vs with taking into account the competing risk of death for different couples (τ, τ'). The dashed lines represent the error made for individual 12 in our hypothetical family.

References

- Elisabeth B Claus, N Risch, and W Douglas Thompson. Genetic analysis of breast cancer in the cancer and steroid hormone study. *American journal of human genetics*, 48(2):232, 1991.
- DF Easton, DT Bishop, D Ford, and GP Crockford. Genetic linkage analysis in familial breast and ovarian cancer: results from 214 families. the breast cancer linkage consortium. *American journal of human genetics*, 52(4):678, 1993.
- INED. Death incidence in France, period 2012-2014. https://www.ined.fr/en/everything_about_population/data/france/deaths-causes-mortality/mortality-tables/, 2017.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Liviu R Totir, Rohan L Fernando, and Joseph Abraham. An efficient algorithm to compute marginal posterior genotype probabilities for every member of a pedigree with loops. *Genetics Selection Evolution*, 41(1):52, 2009.

Acknowledgment

This work has been funded by Institut Curie (M2 internship) and Ligue Nationale Contre le Cancer (PhD grant).
[correspondance: lefebalex@gmail.com or nuel@math.cnrs.fr]